1. Describe all the methods you have implemented. (60%)

   I tried the pretrain model from hfl called "hfl/chinese-bert-wwm-large" and "hfl/chinese-bert-wwm". Although they claim they are better, they actually are not. These two model use BertTokenizer and BertModelforMultipleChoice to load. However, the config of their tokenizer is not complete. To use it, we have to add max_length=512 and truncation=True. Otherwise, an error will occur. Other than that, I also tried pretrain model from CKIP lab called "ckiplab/albert-base-chinese" and "ckiplab/bert-base-chinese". The albert one performs significantly better than bert and pass the baseline.

2. Did you preprocess your data from the dataset? Why? And how?
(Did you encounter the problem that the input length is longer than the maximum sequence length of the model you use? How did you solve this problem?) (30%)

Yes, I preprocess my data by change simple Chinese to traditional Chinese using Opencc library because CKIP lab pretrain its models on traditional Chinese data. Furthermore, I created new json files for huggingface to load, flatten each questions base on the same context. Also, I actually train three different model for different number of choices such as 2,3 and 4.

As I stated in question 1, I added truncation=True and max_length=512 in the tokenizer.

3. What difficulties did you encounter in this assignment? How did you solve it? (10%)
The biggest difficulty I encountered is truncation issue. Since the official example of tokenizer with truncation does not add max_length=512, some of the pretrained models will have problem. To be caution, one has to specify the max_length to enable truncation for every models.